Evolutionary Learning

Changshui Zhang, Jianwen Zhang, Yangqing Jia

State Key Laboratory on Intelligent Technology and Systems Tsinghua National Laboratory for Information Science and Technology (TNList) Department of Automation, Tsinghua University, Beijing 100084, China

November 7, 2009

Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

크

Outline of Part I

A dilemma between learning algorithms and practical systems

Evolutionary learning



Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

크

Outline of Part I

A dilemma between learning algorithms and practical systems

2 Evolutionary learning



Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Outline of Part I

A dilemma between learning algorithms and practical systems

2 Evolutionary learning



Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

(日) (四) (E) (E)

Outline of Part II



5 Evolutionary clustering

Online evolutionary exponential family mixture
 A density estimation viewpoint to clustering
 The roles of data and model: two general approaches
 Experiments

Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

Outline of Part II





Online evolutionary exponential family mixture
 A density estimation viewpoint to clustering
 The roles of data and model: two general approache
 Experiments

Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

Outline of Part II



Pioneer works



Online evolutionary exponential family mixture

- A density estimation viewpoint to clustering
- The roles of data and model: two general approaches
- Experiments

Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

Outline of Part III



8 A brief review of SSL

Semi-supervised evolutionary classification

00 Experiments

Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

Outline of Part III





9 Semi-supervised evolutionary classification

00 Experiments

Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Outline of Part III



Evolutionary classification



Semi-supervised evolutionary classification

Part I: Introduction Part II: Unsupervised Evolutionary Learning Part III: Semi-Supervised Evolutionary Learning

Outline of Part III



Evolutionary classification



Semi-supervised evolutionary classification

10 Experiments

Part I

Introduction

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

æ

Static learning & dynamic systems

• Conventional learning on static data

$$\{X_{train}, X_{test}\} \stackrel{i.i.d.}{\sim} P(x)$$

Practical systems: dynamic

$$\{X_{train}^{1}, X_{test}^{1}\} \stackrel{i.i.d.}{\sim} P_{1}(x)$$
$$\{X_{train}^{2}, X_{test}^{2}\} \stackrel{i.i.d.}{\sim} P_{2}(x)$$

.

$$\{X_{train}^t, X_{test}^t\} \stackrel{i.i.d.}{\sim} P_t(x)$$

 $P_1(x) \rightarrow P_2(x) \rightarrow \ldots \rightarrow P_t(x)$

< 注 > < 注 >



 Topics on Internet. Google Trends: five keywords ("nuclear", "economics", "oil", "tennis", "NBA")



・ロン ・雪 ・ ・ ヨ ・

Blog, social networks, vedio,...

Dilemma

"Static" is the typical setting we think about learning problems.

"Dynamic" is the typical behavior of practical systems implementing learning algorithms.

・ロ・ ・ 四・ ・ ヨ・ ・ 日・ ・

Question:

Is it necessary to introduce a new algorithm especially for dynamic problems?

Possible solvers?

- Applying static algorithms on the whole collection of data.
 Non I.I.D.
- 2. Applying static algorithms on data at each time snapshot $X^t \to \mathcal{M}^t(x)$ or $p(\mathcal{M}^t | X^t)$.
 - Sample size (especially the labeled samples)
 - Time evolving mechanism.
 - Predicting (dynamic predicting rather than the classifier's predicting)

< 日 > < 回 > < 回 > < 回 > < 回 > <

크

Evolutionary learning

Evolutionary learning

The learning problem on time evolving data.

Learning tasks

Static performance

$$\sum_{t} loss(\mathcal{P}_{t}, \mathcal{M}_{t})$$

• Learning time evolving mechanism

$$p\left(\mathcal{P}_{t+1}\middle|\mathcal{P}_t,\ldots,\mathcal{P}_1\right)$$

・ロン ・雪 ・ ・ ヨ ・



Online

$$\mathcal{M}_t | X_t, \ldots, X_1$$

Off-line

$$\mathcal{M}_1,\ldots,\mathcal{M}_t,\ldots,\mathcal{M}_T|X_1,\ldots,X_t,\ldots,X_T$$

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・ ・

- (Semi-)supervised: classification / regression
- Unsupervised: clustering / density estimation

A dilemma between learning algorithms and practical systems Evolutionary learning Analogs Learning on time series



- Incremental learning
- Online learning
- Learning on time series

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Incremental learning Online learning Learning on time series

A B > A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A
 A

Incremental learning

Data are organized into batch or stream mode due to limited memory or on-line setting, etc.

Differences from evolutionary learning

- Assumption: all data comes from an identical distribution.
- Output: a single output (classifier / regressor / data partition / density model /) for the entire data set.
- Not put forward for temporary data, on the contrary, a good incremental algorithm is expected to be insensitive to the order of data.

A dilemma between learning algorithms and practical systems Evolutionary learning Analogs Learning on time series

Online learning I

A special type of practical learning environment: a sequence of consecutive rounds. When predicting at *t*, leaner can only observe samples up to *t*. Leaner dynamically updates, to pursuit a low *regret*.



Difference between online learning and evolutionary learning.

1. The data access of evolutionary learning can be online or off-line.

A B > A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A
 A

Incremental learning Online learning Learning on time series

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Online learning II

- 2. At each snapshot *t*, online learning observes a single sample, while evolutionary learning observe a package of samples.
- 3. Online learning requires a feedback for each sample, while in evolutionary learning, the package of samples may be labeled, or unlabeled, or partially labeled.
- 4. Online learning minimize a *regret* relative to the corresponding offline learning. No generalization problem exists at each snapshot *t*. Evolutionary learning aims to achieve two objectives: generalized loss at each snapshot *t*; the time evolving mechanism.

Incremental learning Online learning Learning on time series

・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

Learning on time series

- Time series: at *t*, a single observation *x_t*
- Evolutionary learning: at t, a distribution \mathcal{P}_t (i.i.d. samples X_t)



Part II

Unsupervised Evolutionary Learning

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

크

Outline of Part II

Pioneer works

5 Evolutionary clustering

Online evolutionary exponential family mixture
 A density estimation viewpoint to clustering
 The roles of data and model: two general approaches
 Experiments

Outline of Part II





Online evolutionary exponential family mixture
 A density estimation viewpoint to clustering
 The roles of data and model: two general approaches
 Experiments

・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

Outline of Part II





Online evolutionary exponential family mixture

- A density estimation viewpoint to clustering
- The roles of data and model: two general approaches

Experiments

Pioneer works on evolutionary learning

- Dynamic topic model [BL06] (ICML'06)
- Segmentation of image sequences [GG06] (TPAMI'06)
- Evolutionary clustering [CKT06, CSZ⁺07, AX08] (KDD'06, KDD'07, SDM'08).

<ロ> <同> <同> <同> < 同> < 同> 、

Dynamic topic model [BL06] (ICML'06)

A dynamic extension to the topic model *Latent Dirichlet Allocation* [BNJL03] (JMLR'03)



・ロ・ ・ 四・ ・ ヨ・ ・ 日・ ・

Segmentation of image sequence [GG06] (TPAMI'06)

Each frame image is modeled by a Bayesian mixture model.

Bayesian Gaussian mixture model

$$z \sim Multnomial(\alpha_1, \dots, \alpha_K), \quad x \sim \mathcal{N}(\mu_z, \Sigma_z)$$
$$\alpha \sim Dirichlet(\gamma), \quad \left(\mu_z, \Sigma_z^{-1}\right) \sim \mathcal{NW}(\Theta_z)$$

where $\mathcal{NW}(\Theta_z)$ is the *Normal Wishart* distribution (the conjugate prior for Gaussian).

[GG06] let $\alpha_t \sim \text{Dirichlet}(\gamma_{t-1}), \quad \left(\mu_{z,t}, \Sigma_{z,t}^{-1}\right) \sim \mathcal{NW}(\Theta_{z,t-1}).$ Notice: in both [BL06] and [GG06], component number is fixed along time.

Algorithms

Evolutionary Clustering

Arising from data mining [CKT06] (KDD'06).

The target of evolutionary clustering

- For each time epoch t, output a partition Π_t on X_t ;
- A trade-off between good partition quality and preserving the smoothness of {Π_t}^T_{t=1} along t.

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Cluster tracking.

Algorithms

Why is evolutionary clustering meaningful?

- Interpretability of machine learning and data mining algorithms.
- The importance of consistency between clustering results along time, especially for visualization.
- Clustering tracking provides powerful approach to dynamic network behavior analysis.

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Algorithms

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Why is evolutionary clustering challenging?

- The variation of data size.
- The variation of cluster number.
- Dynamic behavior of clusters along time, e.g., birth, merging, death, etc.

Algorithms

Evolutionary clustering of [CKT06]

An abstract framework

$$\max_{C_t} sq(C_t, M_t) - \lambda \cdot hc(C_{t-1}, C_t).$$

C_t: partition at *t*; *M_t*: similarity matrix between samples; *sq*: *snapshot quality*; *hc*: *historical cost*.

Evolutionary k-means

$$\min \sum_{i=1}^{n_t} \|x_i - m_{c(x_i)}\| + \lambda \sum_{c=1}^{C_t} \|m_c^t - m_{c'}^{t-1}\|$$
$$c' = \arg \min_{z} \|m_c^t - m_{z'}^{t-1}\|$$

Changshui Zhang, Jianwen Zhang, Yangqing Jia

Evolutionary Learning

Algorithms

Evolutionary spectral clustering [CSZ+07]

• Original Normalized Cut

(

$$Cost_{NC} = k - \operatorname{Tr}\left[Z^{\top}\left(D^{-\frac{1}{2}}WD^{-\frac{1}{2}}\right)Z\right]$$

Preserving Cluster Quality (PCQ)

$$Cost_{NC} = \alpha \cdot NC_t |_{Z_t} + \beta \cdot NC_{t-1} |_{Z_t}$$
$$= k - Tr \left[Z_t^\top \left(\alpha D_t^{-\frac{1}{2}} W_t D_t^{-\frac{1}{2}} + \beta D_{t-1}^{-\frac{1}{2}} W_{t-1} D_{t-1}^{-\frac{1}{2}} \right) Z_t \right]$$

Preserving Cluster Membership (PCM)

$$Cost_{NC} = \alpha \cdot NC_t |_{Z_t} + \beta \cdot CT (Z_t, Z_{t-1})$$
$$= k - Tr \left[Z_t^\top \left(\alpha D_t^{-\frac{1}{2}} W_t D_t^{-\frac{1}{2}} + \beta Z_{t-1} Z_{t-1}^\top \right) Z_t \right]$$

Algorithms

Two different settings about evolutionary clustering

[CSZ⁺07] differs with [CKT06] in the basic data setting:

 Evolutionary spectral clustering of [CSZ⁺07]



 Evolutionary clustering of [CKT06]



・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・
Algorithms

Problems not made clear

- What's evolving?
- What dose "smoothness" mean?
- Evolutionary k-means, evolutionary GMM, etc. Is there a general approach?

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

- The behavior of clusters
- The variation of cluster number and data size

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

Online evolutionary exponential family mixture(Jianwen Zhang, etc. IJCAI, 09)

- A density estimation viewpoint to clustering problem
- Online evolutionary exponential family mixture
- Experiments

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

크

Density estimation via mixture models

Density estimation (Fisher-Wald setting)

$$\min_{\Xi} \mathcal{L}(\Xi) = -\int \log p(\mathbf{x}; \Xi) dF(\mathbf{x})$$
$$= KL(f||p) - \int f(x) \log f(x) dx$$

Mixture models

$$p(\mathbf{x}; \mathbf{\Xi}) = \sum_{z}^{C} \alpha_{z} p(\mathbf{x}; \theta_{z}), \text{ with } \sum_{z} \alpha_{z} = 1$$

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Solver: Expectation-maximization (EM)

A varational upper bound

 $\mathcal{L}(\Xi) \leq \mathcal{G}\left(q_x(\cdot), \Xi\right)$

where $q_x(\cdot)$ is a probability distribution for *z*.

$$\begin{array}{l} \textbf{E}\text{-step: } q_{\textbf{x}}^{[t+1]}(\cdot) \leftarrow \operatorname*{arg\,min}_{q_{\textbf{x}}(\cdot)} \mathcal{G}(q_{\textbf{x}}(\cdot), \Xi^{[t]}) \\ \textbf{M}\text{-step: } \Xi^{[t+1]} \leftarrow \operatorname*{arg\,min}_{\Xi} \mathcal{E}(q_{\textbf{x}}^{[t+1]}(\cdot), \Xi) \end{array}$$

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

E-step: as a clustering problem

• If no additional constraints on $q_x(\cdot)$,

$$q_{\mathbf{x}}^{[t+1]}(z) = p(z|\mathbf{x}, \Xi^{[t]})$$

This case is called "soft-clustering";

• If we constraint that $q_x(z) \in \{0, 1\}$,

$$q_{\mathbf{x}}^{[t+1]}(z) = \mathbf{I}_{\left[z = \arg\max_{z} \rho(z | \mathbf{x}; \Xi^{[t+1]})\right]}$$

This case is called "hard-clustering";

A density estimation viewpoint to clustering Pioneer works Evolutionary clustering Online evolutionary exponential family mixture Pioneer works Evolutionary clustering Historical data dependent (HDD) Historical model dependent (HMD) Experiments

M-step

 Solution with closed form when using *exponential family mixture* (EFM):

$$\mu_z^{[t+1]} = \nabla \Psi(\theta_z^{[t+1]}) = \frac{\mathsf{E}_f \left[q_{\mathsf{x}}(z) T(\mathsf{x}) \right]}{\mathsf{E}_f \left[q_{\mathsf{x}}(z) \right]}$$
$$\alpha_z^{[t+1]} = \mathsf{E}_f \left[q_{\mathsf{x}}(z) \right]$$

• Exponential family mixture (EFM)

$$p(\mathbf{x}; \Psi, \Xi) = \sum_{z}^{C} \alpha_{z} p_{\Psi}(\mathbf{x}; \theta_{z})$$

 $p_{\Psi}(\mathbf{x}; \theta_z)$ belongs to a same *expoenntial family*:

$$p_{\Psi}(\mathbf{x}; oldsymbol{ heta}) = \exp\left\{ \langle oldsymbol{ heta}, \mathcal{T}(\mathbf{x})
angle - \Psi(oldsymbol{ heta})
ight\} p_0(\mathcal{T}(\mathbf{x})).$$

Gaussian, Multinomial, Poisson, Binomial, Exponential, Dirichlet, ...

Changshui Zhang, Jianwen Zhang, Yangqing Jia Evolutionary Learning

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

Two general approaches

- Data & model: Static: f(x), p(x) Evolutionary: f⁽¹⁾(x), f⁽²⁾(x), p⁽¹⁾(x), p⁽²⁾(x)
- Two possible approaches

$$\mathcal{L} = (1 - \lambda) \cdot \textit{dist} \left(f^{(2)}, p^{(2)} \right) + \frac{\lambda \cdot \textit{dist} \left(f^{(1)}, p^{(2)} \right)}{\lambda \cdot \textit{dist} \left(p^{(1)}, p^{(2)} \right)}$$

Historical data dependent (HDD) Historical model dependent (HMD)





Two possible approaches

・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

・ロ > ・ 一部 > ・ モ = > ・ ・ モ > ・

Historical data dependent (HDD)

Loss function

$$\begin{split} \mathcal{L}_{d} &= (1 - \lambda) \cdot \textit{KL}(f^{(2)} || \textit{p}^{(2)}) + \lambda \cdot \textit{KL}(f^{(1)} || \textit{p}^{(2)}) \\ &= -\mathbf{E}_{\tilde{f}_{\lambda}}[\log \textit{p}^{(2)}(\mathbf{x}; \Xi^{(2)})] \end{split}$$

where $\tilde{f}_{\lambda}(\mathbf{x}) = (1 - \lambda)f^{(2)}(\mathbf{x}) + \lambda f^{(1)}(\mathbf{x}).$

- Meaning
 Using EFM p⁽²⁾ to estimate *f*_λ. (Recall the static clustering
 as using EFM p to estimate f: L = -E_f[log p(x; Ξ)])
- Solver

EM.

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

・ロト ・ 日 ・ ・ 目 ・ ・ 日 ・ ・

æ

Solution to HDD

- E-step: the same as clustering using EFM.
- M-step

$$\begin{aligned} \alpha_z^{(2),[t+1]} &= \mathbf{E}_{\tilde{f}_{\lambda}} \left[q_{\mathbf{x}}^{[t+1]}(z) \right] \\ \mu_z^{(2),[t+1]} &= \frac{\mathbf{E}_{\tilde{f}_{\lambda}} \left[q_{\mathbf{x}}^{[t+1]}(z) T(\mathbf{x}) \right]}{\mathbf{E}_{\tilde{f}_{\lambda}} \left[q_{\mathbf{x}}^{[t+1]}(z) \right]} \end{aligned}$$

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

크

Historical model dependent (HMD)

Loss function

$$\mathcal{L}_m = (1 - \lambda) \mathcal{KL}(f^{(2)}, \boldsymbol{p}^{(2)}) + \lambda \boldsymbol{d}_{\mathsf{EMD}}(\boldsymbol{p}^{(1)}, \boldsymbol{p}^{(2)}).$$

where

$$d_{\text{EMD}}(p^{(1)}, p^{(2)}) = \min_{\mathbf{w}} \sum_{l, z} w_{lz} d(p(\mathbf{x}; \theta_l^{(1)}), p(\mathbf{x}; \theta_z^{(2)}))$$

s.t. $w_{lz} \ge 0, \sum_{z} w_{lz} = \alpha_l^{(1)}, \sum_{l} w_{lz} = \alpha_z^{(2)}$

The equivalent problem:

$$\min_{\Xi^{(2)}, \mathbf{w}} \mathcal{L}'_m(\Xi^{(2)}, \mathbf{w}) = (1 - \lambda) \mathcal{K}L(f^{(2)}(x), p^{(2)}(x; \Xi^{(2)})) + \lambda \sum_{l, z} w_{lz} \mathcal{K}L(p(x; \theta_l^{(1)}), p(x; \theta_z^{(2)}))$$

s.t. . . .

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

Optimization

• Variational upper bound by introducing $q_x(z)$:

$$\mathcal{L}'_{m}(\mathbf{w}, \mathbf{\Xi}^{(2)}) \leq \mathcal{G}\left(q_{\mathbf{x}}(\cdot), \mathbf{\Xi}^{(2)}, \mathbf{w}
ight)$$

- Alternatively optimize w.r.t. $q_x(z)$, $\Xi^{(2)}$ and **w**.
- q-step: posterior
- w-step: min_w $\sum_{l,z} w_{lz} \mathcal{KL}\left(p(\mathbf{x}; \theta_l^{(1)}), p(\mathbf{x}; \theta_z^{(2)})\right)$ s.t. ...
- Ξ-step:

$$\alpha_{z}^{(2),[t+1]} = \mathbf{E}_{f^{(2)}}[q_{\mathbf{x}}^{[t+1]}(z)]$$
$$\mu_{z}^{(2),[t+1]} = \frac{(1-\lambda)\mathbf{E}_{f^{(2)}}[q_{\mathbf{x}}^{[t+1]}(z)T(\mathbf{x})] + \lambda \sum_{l} w_{lz}^{[t+1]}\mu_{l}^{(1)}}{(1-\lambda)\mathbf{E}_{f^{(2)}}[q_{\mathbf{x}}^{[t+1]}(z)] + \lambda \sum_{l} w_{lz}^{[t+1]}}$$

A density estimation viewpoint to clustering Pioneer works Evolutionary clustering Online evolutionary exponential family mixture Experiments

• Data: NSF research awards abstracts.

- 13 years, D = 19728, W = 15412
- Features
 - Word count (For Mixture of multinomial)
 - tf-idf (For k-means)
- Models
 - *k-means* → Evolutionary k-means
 - Multinomial mixture model \rightarrow Evolutionary MMM

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

	A density estimation viewpoint to clustering
Pioneer works	The roles of data and model: two general approaches
Evolutionary clustering	Historical data dependent (HDD)
Online evolutionary exponential family mixture	Historical model dependent (HMD)
	Experiments

Evaluation

Normalized mutual information (NMI) at each epoch

$$NMI = \frac{I(\Upsilon; \hat{\Upsilon})}{\sqrt{H(\Upsilon) \cdot H(\hat{\Upsilon})}} \approx \frac{\sum_{h,c} \log\left(\frac{n \cdot n_{h,c}}{n_{h} n_{c}}\right)}{\sqrt{\left(\sum_{h} n_{h} \log \frac{n_{h}}{n}\right) \left(\sum_{c} n_{c} \log \frac{n_{c}}{n}\right)}}$$

Data measured temporal loss (DTL)

$$DTL = -\mathbf{E}_{f^{(1)}} \left[\log p^{(2)(\mathbf{x})} \right] = KL \left(f^{(1)} || p^{(2)} \right) - Const$$

Model measured temporal loss (MTL)

$$\textit{MTL} = \textit{d}_{\mathsf{EMD}}\left(\textit{p}^{(1)},\textit{p}^{(2)}
ight)$$

・ロ・ ・ 四・ ・ ヨ・ ・ 日・ ・

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

▲ 同 ▶ ▲ 三 ▶

Results: Evolutionary k-means

Snapshot performance: NMI



 Pioneer works
 A density estimation viewpoint to clustering

 Evolutionary clustering
 The roles of data and model: two general approaches

 Historical data dependent (HDD)
 Historical model dependent (HMD)

 Experiments
 Experiments

Temporal performance: DTL & MTL



・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

æ

Pioneer works Evolutionary clustering Online evolutionary exponential family mixture Online evolutionary exponential family mixture

Results: Evolutionary MMM

Snapshot performance: NMI



A (1) > A (2) > A

ъ

 Pioneer works
 A density estimation viewpoint to clustering

 Evolutionary clustering
 The roles of data and model: two general approaches

 Historical data dependent (HDD)
 Historical model dependent (HMD)

 Experiments
 Experiments

Temporal performance: DTL & MTL





・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

æ

A density estimation viewpoint to clustering The roles of data and model: two general approaches Historical data dependent (HDD) Historical model dependent (HMD) Experiments

・ロト ・ 一 ト ・ ヨ ト ・ ヨ ト



- A density estimation view point to clustering and evolutionary clustering problem
- Two general approaches based on exponential family mixture models
- Validated by different EFMs on real text data
- Unsolved problem: tracking of clusters

Part III

Supervised Evolutionary Learning

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

æ

Changshui Zhang, Jianwen Zhang, Yangqing Jia Evolutionary Learning

Outline of Part III



8 A brief review of SSL

Semi-supervised evolutionary classification

00 Experiments

Changshui Zhang, Jianwen Zhang, Yangqing Jia Evolutionary Learning

・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

Outline of Part III





A brief review of SSL

Changshui Zhang, Jianwen Zhang, Yangqing Jia **Evolutionary Learning**

・ロト ・ 聞 ト ・ ヨ ト ・ ヨ ト

Outline of Part III





Semi-supervised evolutionary classification

10 Experiments

Changshui Zhang, Jianwen Zhang, Yangqing Jia Evolutionary Learning

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Outline of Part III





Semi-supervised evolutionary classification

10 Experiments

Changshui Zhang, Jianwen Zhang, Yangqing Jia Evolutionary Learning

A B A B A B A B A
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Evolutionary classification

Evolutionary classification

Evolutionary classification learns a chain of evolving classifiers for different time periods.

Why necessary?

- Using historical classifiers makes little sense:
 - \rightarrow i.i.d. assumption dose not hold.
- Using all historic data also makes little sense:
 - \rightarrow i.i.d. assumption dose not hold.
- Historical classification information may help:
 - \rightarrow distribution evolves slowly (assumed but make sense)
- Historical labels may be utilized rather than discarded

A Brief Review of SSL

- Data: $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_l, \mathbf{x}_{l+1}, \cdots, \mathbf{x}_n\}, \mathcal{Y} = \{y_i\}_{i=1}^l$.
- Require: $\{y_i\}_{i=l+1}^n$, or $f(\mathbf{x})$.
- Representative Work: Manifold Regularization [BNS06]:

$$\min_{f} \mathcal{J}(f) = \sum_{i=1}^{I} \mathcal{L}(\mathbf{y}_i, \mathbf{x}_i, f) + \gamma_{\mathcal{A}} \|f\|_{\mathcal{K}}^2 + \gamma_{I} \|f\|_{I}^2 \quad (1)$$

where

$$\|f\|_{I}^{2} = \frac{1}{n^{2}} \sum_{i,j=1}^{n} (f(\mathbf{x}_{i}) - f(\mathbf{x}_{j}))^{2} W_{ij} = \frac{1}{n^{2}} \mathbf{f}^{\top} \mathbf{L} \mathbf{f}$$
(2)

・ロ > ・ 一部 > ・ モ = > ・ ・ モ > ・

크

is the spatial regularizer calculated via graph Laplacian.

Semi-supervised Evolutionary Classification: Basic Settings(Yangqing Jia. etc. IJCAI, 09)

- In general, we aim to learn a function *F*(*t*) : ℝ → *H_K*, where *H_K* is the Hilbert space of functions *f* : ℝ^d → ℝ associated with a pre-defined kernel *K*.
- Specifically, *F*(*t*) = *f_t* gives the classification function for each time *t*.
- Data: $\mathcal{X} = \{\mathcal{X}_1, \mathcal{X}_2, \cdots, \mathcal{X}_T\}$ from *T* discrete "frames", where each $\mathcal{X}_t = \{\mathbf{x}_{1,t}, \mathbf{x}_{2,t}, \cdots, \mathbf{x}_{n_t,t}\}$, with the first I_t labeled as $\mathcal{Y}_t = \{y_i\}_{i=1}^{l_t}$.
- When the time *t* takes discrete values, the goal is to find *T* classification functions $\{\mathcal{F}(t)\}_{t=1}^{T} = \{f_t(\mathbf{x})\}_{t=1}^{T}$.

イロト イヨト イヨト

Evolutionary Smoothness Assumption

Similar to the smoothness assumption in the general learning problems, we extend it to the evolutionary data.

Smoothness Assumption

Two data points are likely to have similar labels if they are close to each other.

Evolutionary Smoothness Assumption

Two classification functions f_{t_1} and f_{t_2} are likely to be similar if the times t_1 and t_2 are close.

A B > A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

Temporal Regularizer

 A direct approach to carry out the assumption above is to use the integral of ||∂F/∂t||² over time t:

$$\int_{1}^{T} \left\| \frac{\partial \mathcal{F}}{\partial t} \right\|^{2} \mathrm{d}t \tag{3}$$

• When *t* is discrete, we use the backward difference to approximate the above integral as

$$\int_{1}^{T} \left\| \frac{\partial \mathcal{F}}{\partial t} \right\|^{2} \mathrm{d}t \approx \frac{1}{T-1} \sum_{t=2}^{T} \|f_{t}(\cdot) - f_{t-1}(\cdot)\|_{K}^{2} \qquad (4)$$

This serves as a temporal regularizer for the learning algorithm.

Offline Algorithm

 Taking into consideration the evolutionary smoothness assumption, we propose an offline learning algorithm simultaneously learns the *T* classification functions by minimizing the following objective function w.r.t. {*f_t*}^{*T*}_{*t*=1}:

$$\mathcal{J}_{\text{off}} = \sum_{t=1}^{T} \left[\sum_{i=1}^{n_t} \mathcal{L}(y_{i,t}, \mathbf{x}_{i,t}, f_t(\cdot)) + \gamma_I \| f_t(\cdot) \|_I^2 \right] \\ + \gamma_A \| f_1 \|_K^2 + \gamma_A \sum_{t=2}^{T} \| f_t(\cdot) - f_{t-1}(\cdot) \|_K^2$$
(5)

 It is not difficult to observe that the problem is convex if the loss function *L* is convex, thus the global optimal solution can be found.

Online Algorithm

- As the data accumulates, the learning problem scales up fast and renders the offline algorithm impractical.
- Thus, we propose an online algorithm that learns one classifier for the current frame with the historic information fixed. For time *t*, the algorithm minimizes the following objective function w.r.t. *f_t*:

$$\mathcal{J}_{on}^{(t)} = \sum_{i=1}^{n_t} \mathcal{L}(\mathbf{y}_{i,t}, \mathbf{x}_{i,t}, f_t(\cdot)) + \gamma_I \| f_t(\cdot) \|_I^2 + \gamma_A \| f_t(\cdot) - f_{t-1}(\cdot) \|_K^2$$
(6)

A B > A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

The Representer Theorem

It can be proved that in the evolutionary case, the representer theorem still holds, enabling us to convert the problem from finding the functions $\{f_t\}_{t=1}^T$ to simply finding a set of expansion coefficients.

Representer Theorem

Assume that the regularizers $||f_t(\cdot)||_l^2$ $(1 \le t \le T)$ are carried out empirically using the graph Laplacian as $\mathbf{f}_t^\top \mathbf{L}_t \mathbf{f}_t$, then for each frame *t*, the minimizer of (6) admits an expansion

$$f_t^*(\mathbf{x}) = \sum_{k=1}^t \sum_{i=1}^{n_t} \alpha_{i,k} K(\mathbf{x}, \mathbf{x}_{i,k}), \quad \forall 1 \le t \le T$$
(7)

given a predefined Mercer kernel $K(\cdot, \cdot)$.

Solution Space Shrinking

- Problem: as the data accumulates, the scale of the expansion subspace $S \subset \mathcal{H}_K$ that is constructed by the kernel functions $\{K(\cdot, \mathbf{x}_{i,k})\}_{i=1}^{n_k} \underset{k=1}{t}$ may grow too large for learning.
- Thus, we manually impose a *representer constraint* to the optimization problem by shrinking the solution space:

$$f^*(\cdot; t) = \underset{f_t(\cdot) \in \mathcal{S}_t}{\operatorname{arg\,min}} \quad \mathcal{J}_{\operatorname{on}}^{(t)} \tag{8}$$

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

where $S_t = \text{span}\{K(\cdot, x_{i,t}) | x_{i,t} \in \mathcal{X}_t\}$ is spanned by the kernel functions corresponding to the dataset \mathcal{X}_t

Kernel Based Representation

•
$$f_t(\cdot) = \sum_{i=1}^{n_t} \alpha_{i,t} K(\cdot, \mathbf{x}_{i,t})$$

The objective function can be written as

$$\mathcal{J}_{\text{on}}^{(t)} = (\mathbf{K}_{t} \boldsymbol{\alpha}_{t} - \mathbf{y}_{t})^{\top} \mathbf{C}_{t} (\mathbf{K}_{t} \boldsymbol{\alpha}_{t} - \mathbf{y}_{t}) + \gamma_{A} \boldsymbol{\alpha}_{t}^{\top} \mathbf{K}_{t} \boldsymbol{\alpha}_{t} - 2 \gamma_{A} \boldsymbol{\alpha}_{t}^{\top} \mathbf{K}_{t,t-1} \boldsymbol{\alpha}_{t-1} + \frac{\gamma_{I}}{n_{t}^{2}} \boldsymbol{\alpha}_{t}^{\top} \mathbf{K}_{t}^{\top} \mathbf{L}_{t} \mathbf{K}_{t} \boldsymbol{\alpha}_{t}$$
(9)

And the solution is

$$\boldsymbol{\alpha}_{t}^{*} = (\mathbf{K}_{t}^{\top}(\mathbf{C}_{t} + \frac{\gamma_{I}}{n_{t}^{2}}\mathbf{L}_{t})\mathbf{K}_{t} + \gamma_{A}\mathbf{K}_{t})^{-1} \\ (\mathbf{K}_{t}\mathbf{C}_{t}\mathbf{y}_{t} + \gamma_{A}\mathbf{K}_{t,t-1}\boldsymbol{\alpha}_{t-1})$$
(10)

・ロ > ・ 一部 > ・ モ = > ・ ・ モ > ・

æ

Toy Data



Figure: Standard SSL on each frame



Figure: SSL-E

▲□ → ▲ □ → ▲ □ →

æ

Real-world Dataset

- We crawled posts from 5 online mailing lists dated from 2003.9 to 2008.9, each month as a time frame.
- There are 161,675 posts (data points) for the learning task.
- We applied two other methods to compare with our algorithm:
 - classical SSL algorithm performed on each frame;
 - a naive evolutionary SSL extension that simply uses the data from the current frame and its preceding frame.

Label	Date	Total Posts	Posts/Month	
audiophiles	2005.6 - 2008.9	35,666	892	
listening-l	2003.7 - 2008.9	25,738	422	
sqlite-users	2003.10 - 2008.9	32,655	544	
tutor-python	2004.12 - 2008.9	26,314	572	
wine-devel	2003.9 - 2008.9	41,302	677	

A B > A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A

Experimental Results: Statistics

Table: Experimental results on the evolutionary mailing list data using SSL, SSL-naive, and SSL-E.

class	aud/lis	aud/sql	aud/tut	aud/win
SSL	97.03(4.42)	96.12(5.19)	94.04(7.05)	93.14 (8.20)
SSL-naive	97.77(2.98)	95.99(6.39)	95.47(6.33)	93.18 (10.28)
SSL-E	98.98 (2.03)	98.80 (2.97)	97.95(4.87)	97.11 (6.89)
class	lis/sql	lis/tut	lis/win	sql/tut
SSL	84.40 (9.02)	81.50(9.89)	78.41(10.52)	63.00(5.26)
SSL-naive	86.30(7.39)	85.77 (8.38)	81.19(9.86)	64.52(3.41)
SSL-E	92.74 (6.41)	91.30 (6.65)	91.41 (6.89)	76.87 (4.10)
class	sql/win	tut/win	Multi-class	
SSL	64.42(4.97)	63.33(4.40)	77.49 (5.02)	
SSL-naive	68.45(4.16)	65.57(3.51)	78.13 (4.89)	
SSL-E	80.87 (8.16)	78.56 (7.13)	86.10 (3.21)	

Changshui Zhang, Jianwen Zhang, Yangqing Jia

Evolutionary Learning

・ロ > ・ 一部 > ・ モ = > ・ ・ モ > ・
Evolutionary classification A brief review of SSL Semi-supervised evolutionary classification Experiments

Experimental Results: Spatial and Temporal Loss



Figure: The spatial cost, temporal cost and AUC value vs. time on the *sql/win* data.

A (1) > A (2) > A

Changshui Zhang, Jianwen Zhang, Yangqing Jia Evolutionary Learning

Evolutionary classification A brief review of SSL Semi-supervised evolutionary classification Experiments



- Learning the natural evolutionary information of the data is a new challenge in the machine learning research:
 - the concept drifts during the time and makes a single aggregated classifier inaccurate for long-term prediction;
 - the drifting is smooth in a short time period.
- We proposed a new semi-supervised algorithm for learning a series of evolving classification functions for evolutionary data.
- The proposed algorithm provides much better performances on real-world application in both stability and accuracy.

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Conclusion

Part IV

Conclusion

<ロ> <同> <同> < 同> < 同> < 同> <

3

Changshui Zhang, Jianwen Zhang, Yangqing Jia Evolutionary Learning



- A new type of learning problems goes beyond the I.I.D. assumption. We call it *Evolutionary Learning*.
- A semi-supervised classification algorithm.
- An unsupervised algorithm for clustering / density estimation.
- Just exploratory works. No theoretical analysis provided.

A B > A B > A B >
A
B >
A
B >
A
B >
A
B >
A
B >
A
B >
A
B >
A
B >
A
B >
A
B >
A
B >
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A

Conclusion

References I



Amr Ahmed and Eric Xing.

Dynamic non-parametric mixture models and the recurrent Chinese restaurant process: with applications to evolutionary clustering.

SIAM Conference on Data Minning, 2008.



D.M. Blei and J.D. Lafferty.

Dynamic topic models.

23rd International Conference on Machine Learning (ICML), 2006.



D.M. Blei, A.Y. Ng, M.I. Jordan, and J. Lafferty.

Latent Dirichlet allocation.

Journal of Machine Learning Research, 3(4-5):993-1022, 2003.



Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research*, 7:2434, 2006.



Deepayan Chakrabarti, Ravi Kumar, and Andrew Tomkins.

Evolutionary clustering.

In ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD), pages 554–560. ACM Press New York, NY, USA, 2006.

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・ ・

э

Conclusion

References II



Yun Chi, Xiaodan Song, Dengyong Zhou, Koji Hino, and Belle L. Tseng.

Evolutionary spectral clustering by incorporating temporal smoothness. In ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD). ACM Press New York, NY, USA, 2007.

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

э



J. Goldberger and H. Greenspan.

Context-based segmentation of image sequences.

IEEE Transactions on Pattern Analysis and Machine Intelligence, 28(3):463-468, 2006.